

That which is claimed:

1. A computer implemented method to find a mathematical equation that fits a data set having one dependent variable and at least one independent variable comprising:
 - 5 determining the relative contribution of the at least one independent variable to the dependent variable, and defining separate functions that each describe the contribution of a single independent variable to the dependent variable, wherein the functions used to describe the contribution of an independent variable to the dependent variable are derived using residuals of the dependent variable, wherein the residuals comprise the portion of
 - 10 the dependent variable for which a contributing independent variable has not been defined.
2. The method of claim 1, wherein the analysis of residuals is done sequentially, such that at each stage of the analysis, the residuals comprise contributions from a
 - 15 decreasing number of independent variables.
3. The method of claim 1, wherein the method is automatic in that once a user initiates the analysis by inputting a signal to the computer processor, the processor performs the method with no further input from the user.
 - 20
4. The method of claim 1, further comprising calculating a value for missing data for at least one independent variable.
5. The method of claim 1, further comprising providing a quantitative evaluation of
 - 25 the significance of each independent variable to the equation.
6. A computer implemented method to find a mathematical equation that fits a data set having one dependent variable and at least one independent variable comprising the steps of:
 - 30 (a) identifying the independent variable that makes the largest contribution to the dependent variable as the first most important independent variable;

- (b) plotting the dependent variable versus transformations of the first most important independent variable to determine a function that provides a model having the best fit to the data;
 - (c) identifying the independent variable that makes the next largest contribution to the dependent variable as the next most important independent variable;
 - 5 (d) plotting the residuals of the dependent variable versus transformations of the next most important variable to determine a function that comprises the best fit of the next most important independent variable to the residuals, wherein the residuals of the dependent variable comprise the portion of the dependent variable for which a contributing independent variable has not yet been defined; and
 - 10 (e) repeating steps (c) and (d) to identify the next most important independent variable until an optimal number of independent variables having associated functions to describe the dependent variable have been determined.
- 15 7. The method of claim 6, wherein functions to fit independent variables to the dependent variable or residuals of the dependent variable are chosen from at least one predetermined set of functions.
8. The method of claim 6, wherein step (a) comprises the substeps of:
- 20 (i) plotting the dependent variable versus transformations of each independent variable from the data set;
 - (ii) determining the fit for each independent variable with each of the functions tested in step (i); and
 - (iii) identifying the most important independent variable as the variable having the best fit with at least one of the tested functions.
- 25 9. The method of claim 8, wherein the set of functions used to identify independent variables is smaller than the set of functions used to fit the independent variables to the dependent variable or residuals of the dependent variable.
- 30 10. The method of claim 6, wherein step (c) comprises the substeps of:

- (i) plotting residual values for the dependent variable versus any independent variables that have not been fit to the dependent variable or residuals of the dependent variable;
- (ii) determining the fit for the residual values for the dependent variable 5 versus each of the remaining independent variables; and
- (iii) identifying the next most important independent variable as the variable having the best fit with the residual values for the dependent variable.

11. The method of claim 6, further comprising generating a report comprising at least 10 one equation that includes at least one optimized function for at least one independent variable to describe the value of the dependent variable for the entire data set.

12. The method of claim 11, wherein the report includes generating a list of optimized 15 functions to explain the data set, wherein each of the functions in the list are rated using a predetermined statistical function.

13. The method of claim 12, wherein the list includes functions that include an increasing number of independent variables.

20 14. The method of claim 6, further comprising calculating a value for missing data for at least one independent variable.

15. The method of claim 14, wherein values for missing data are calculated by 25 generating a model or best function without missing the data, and then using the model or best function to derive an approximated value for the missing data.

16. The method of claim 14, wherein the values for missing data are calculated by 30 plotting the independent variable for which the data is missing versus the dependent variable and each of the other independent variables, and estimating a value for the missing data point based on the plot having the best fit.

17. The method of claim 14, wherein the approximated values determined for missing data at one step are used to derive best fit models in subsequent curve-fitting steps.
18. A computer implemented method to find a mathematical equation that fits a data set while minimizing the number of terms in the final model comprising the steps of:
- (a) organizing the data as one dependent variable (y) and at least one independent variable ($x_1, x_2, \dots, x_{n-1}, x_n$);
 - (b) determining which independent variable comprises the most significant contribution to the dependent variable by using a program code that performs the following substeps:
 - (i) plotting the values of the dependent variable against an initial set of selected functions ($F_{initial}$) of each independent variable ($x_1, x_2, x_3, \dots, x_{n-1}, x_n$);
 - (ii) analyzing how well each function describes the values for the dependent variable (y) for each independent variable; and
 - (iii) choosing an independent variable (x_1) which comprises best fit for any one of the predetermined number of analyzed functions;
 - (c) determining a function, $f(x_1)$, and constants, m_1 and b_1 , from an expanded set of functions, which best describes the independent variable comprising the most significant contribution to the dependent variable;
 - (d) determining the residuals $(y - \hat{y}_1)$, where $\hat{y}_1 = m_1 * f(x_1) + b_1$ is the calculated value of (y) for x_1 ;
 - (e) determining the next most significant independent variable by plotting the value of the residuals $(y - \hat{y}_1)$ against an initial set of functions of the remaining independent variables ($x_2, x_3, \dots, x_{n-1}, x_n$) and choosing the independent variable (x_2) which comprises best fit for any one of the predetermined number of analyzed functions;
 - (f) determining a function, $f(x_2)$, and constants, m_2 and b_2 , from an expanded set of functions, which best describes the independent variable comprising the next most significant contribution to the residuals for the dependent variable $(y - \hat{y}_1)$;
 - (g) determining the residuals $(y - \hat{y}_{1,2}) = y - ((m_1 * f(x_1)) + (m_2 * f(x_2)) + b')$;

- (h) plotting selected functions of the remaining independent variables (x_3, \dots, x_{n-1}, x_n) versus the second level residuals ($y - \hat{y}_{1,2}$) in order to determine the next most significant independent variable (x_3);
- 5 (i) determining a function $f(x_3)$, and new constants, m_3 and b_3 , which best describes the mathematical relationship between x_3 and $(y - \hat{y}_{1,2})$ from a second expanded set of pre-selected functions (F_{S2});
- (j) repeating steps (g)-(i) using increasing levels of residuals ($y - y_{1,2,3, \dots, n-1}$) to characterize additional independent variables (x_4, \dots, x_{n-1}, x_n) until an optimal number of functions to describe the dependent variable identified (y) have been and described; and
- 10 (k) generating an equation which includes at least one optimized function for at least one independent variable to describe the value of the dependent variable for the entire data set.
19. A computer-readable medium on which is encoded programming code to find a mathematical equation that fits a data set having one dependent variable and at least one independent variable comprising program code for determining the relative contribution of the at least one independent variable to the dependent variable and for defining separate functions that each describe the contribution of a single independent variable to the dependent variable, wherein the functions used to describe the contribution of an independent variable to the dependent variable are derived using residuals of the dependent variable, wherein the residuals comprise the portion of the dependent variable for which a contributing independent variable has not been defined.
20. The computer-readable medium of claim 19, wherein the program code defines that the analysis of residuals is done sequentially, such that at each stage of the analysis, the residuals comprise contributions from a decreasing number of independent variables.
- 25
21. The computer-readable medium of claim 19, wherein the program code defines that the method is automatic in that once a user initiates the analysis by inputting a signal to the computer processor, the processor performs the method with no further input from the user.
- 30

22. The computer-readable medium of claim 19, further comprising program code for calculating a value for missing data for at least one independent variable.
- 5 23. The computer-readable medium of claim 19, further comprising program code for providing a quantitative evaluation of the significance of each independent variable to the equation.
- 10 24. A computer-readable medium on which is encoded programming code to find a mathematical equation that fits a data set having one dependent variable and at least one independent variable comprising:
- 15 (a) program code for identifying the independent variable that makes the largest contribution to the dependent variable as the first most important independent variable;
- (b) program code for plotting the dependent variable versus transformations of the first most important independent variable to determine a function that provides a model having the best fit to the data;
- 20 (c) program code for identifying the independent variable that makes the next largest contribution to the dependent variable as the next most important independent variable;
- (d) program code for plotting the residuals of the dependent variable versus transformations of the next most important variable to determine a function that comprises the best fit of the next most important independent variable to the residuals, wherein the residuals of the dependent variable comprise the portion of the dependent variable for which a contributing independent variable has not yet been defined; and
- 25 (e) program code for repeating steps (c) and (d) to identify the next most important independent variable until an optimal number of independent variables having associated functions to describe the dependent variable have been determined.

25. The computer-readable medium of claim 24, further comprising program code for choosing functions to fit independent variables to the dependent variable or to residuals of the dependent variable from at least one predetermined set of functions.

5 26. The computer-readable medium of claim 24, wherein the program code for (a) further comprises:

(i) program code for plotting the dependent variable versus transformations of each independent variable from the data set;

10 (ii) program code for determining the fit for each independent variable with each of the functions tested in step (i); and

(iii) program code for identifying the most important independent variable as the variable having the best fit with at least one of the tested functions.

27. The computer-readable medium of claim 26, wherein the set of functions used to identify independent variables is smaller than the set of functions used to fit the independent variables to the dependent variable or residuals of the dependent variable.

28. The computer-readable medium of claim 24, wherein the program code for (c) further comprises:

20 (i) program code for plotting the residual values for the dependent variable versus any independent variables that have not been fit to the dependent variable, or residuals of the dependent variable;

(ii) program code for determining the fit for the residual values for the dependent variable versus each of the remaining independent variables; and

25 (iii) program code for identifying the next most important independent variable as the variable having the best fit with the residual values for the dependent variable.

29. The computer-readable medium of claim 24, further comprising program code for generating a report comprising at least one equation that includes at least one optimized function for at least one independent variable to describe the value of the dependent variable for the entire data set.

30. The computer-readable medium of claim 29, wherein the report includes generating a list of optimized functions to explain the data set, wherein each of the functions in the list are rated using a predetermined statistical function.

5

31. The computer-readable medium of claim 30, wherein the list includes functions that include an increasing number of independent variables.

32. The computer-readable medium of claim 24, further comprising program code for
10 calculating a value for missing data for at least one independent variable.

33. The computer-readable medium of claim 32, further comprising program code to calculate the values for missing data by generating a model or best function without the missing data, and then using the model or best function to derive an approximated value
15 for the missing data.

34. The computer-readable medium of claim 32, further comprising program code to calculate the values for missing data by plotting the independent variable for which the data is missing versus the dependent variable and each of the other independent variables,
20 and estimating a value for the missing data point based on the plot having the best fit.

35. The computer-readable medium of claim 32, further comprising program code to use the approximated values determined for missing data at one step to derive best fit models in subsequent curve-fitting steps.

25

36. A computer-readable medium on which is encoded programming code to find a mathematical equation that fits a data set while minimizing the number of terms in the final model comprising:

(a) program code for organizing the data as one dependent variable (y) and at
30 least one independent variable ($x_1, x_2, \dots, x_{n-1}, x_n$);

(b) program code for determining which independent variable comprises the most significant contribution to the dependent variable by using a program code that performs the following substeps:

- (i) plotting the values of the dependent variable against an initial set 5 of selected functions ($F_{initial}$) of each independent variable ($x_1, x_2, x_3, \dots x_{n-1}, x_n$);
 - (ii) analyzing how well each function describes the values for (y) for each independent variable; and
 - (iii) choosing an independent variable (x_1) which comprises best fit for any one of the predetermined number of analyzed functions;
- 10 (c) program code for determining a function, $f(x_1)$, and constants, m_1 and b_1 , from an expanded set of functions, which best describes the mathematical relationship between the independent variable comprising the most significant contribution to (y);
 - (d) program code for determining the residuals ($y - \hat{y}_1$), where $\hat{y}_1 = m_1 * f(x_1) + b_1$ is the calculated value of (y) for x_1 ;
- 15 (e) program code for determining the next most significant independent variable (x_2) by plotting the value of the residuals ($y - \hat{y}_1$) against an initial set of functions of the remaining independent variable ($x_2, x_3, \dots x_{n-1}, x_n$) and choosing the independent variable (x_2 for example) which comprises best fit for any one of the predetermined number of analyzed functions;
- 20 (f) program code for determining a function, $f(x_2)$, and constants, m_2 and b_2 , from an expanded set functions, which best describes the mathematical relationship between the independent variable comprising the next most significant contribution to (y);
 - (g) program code for determining the residuals $(y - \hat{y}_{1,2}) = y - ((m_1' * f(x_1)) + 25 (m_2' * f(x_2)) + b')$;
- 25 (h) program code for plotting selected functions of the remaining independent variables ($x_3, \dots x_{n-1}, x_n$) versus the second level residuals ($y - \hat{y}_{1,2}$) in order to determine the next most significant independent variable (x_3);
 - (i) program code for determining a function $f(x_3)$, and new constants, m_3 and 30 b_3 , which best describes the mathematical relationship between x_3 and $(y - \hat{y}_{1,2})$ from a second expanded set of pre-selected functions (F_{S2});

- (j) program code for repeating steps (g)-(i) using increasing levels of residuals ($y - y_{1,2,3, \dots, n-1}$) to characterize additional independent variables (x_4, \dots, x_{n-1}, x_n) until an optimal number of functions to describe the dependent variable (y) have been identified and described; and
- 5 (k) program code for generating an equation which includes at least one optimized function for at least one independent variable to describe the value of the dependent variable for the entire data set.